

# UniGeM: Unifying Data Mixing and Selection via Geometric Exploration and Mining

Anonymous ACL submission

## Abstract

The scaling of Large Language Models (LLMs) is increasingly limited by data quality. Most methods handle data mixing and sample selection separately, which can break the structure in code corpora. We introduce **UniGeM**, a framework that unifies mixing and selection by treating data curation as a *manifold approximation* problem without training proxy models or relying on external reference datasets. UniGeM operates hierarchically: **Macro-Exploration** learns mixing weights with stability-based clustering; **Micro-Mining** filters high-quality instances by their geometric distribution to ensure logical consistency. Validated by training 8B and 16B MoE models on 100B tokens, UniGeM achieves **2.0× data efficiency** over a random baseline and further improves overall performance compared to SOTA methods in reasoning-heavy evaluations and multilingual generalization.

## 1 Introduction

The generalization of Large Language Models (LLMs) has traditionally relied on scaling parameters and data volume (Kaplan et al., 2020). However, recent shifts in scaling laws suggest that data quality now constrains model performance more than raw quantity (Hoffmann et al., 2022; Gadre et al., 2024). As high-quality public corpora are nearing depletion (Villalobos et al., 2022), simply aggregating noisy web data yields diminishing returns (Gunasekar et al., 2023). This motivates a focus on **data efficiency**: finding subsets that give more gain per compute (Li et al., 2023; Sorscher et al., 2022).

Most data efficiency approaches fall into two buckets: domain mixing and instance selection (Xie et al., 2023b; Yu et al., 2024; Liu et al., 2025). Data mixing serves as a **macro-distribution balancing** mechanism (Chen et al., 2023), but its coarse granularity often treats domains as flat distributions and overlooks structure within each do-

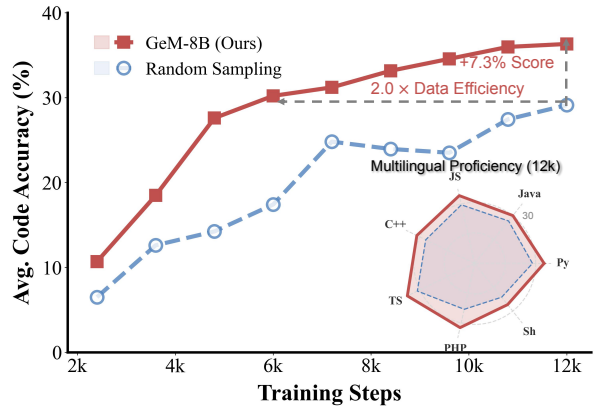


Figure 1: **Average downstream performance:** Random sampling vs. UniGeM for pre-training an 8B MoE model. The inset radar chart shows that UniGeM achieves stronger multilingual performance across programming languages.

main (Diao et al., 2025). Meanwhile, many selection methods require training proxy models or reference datasets, which adds substantial overhead and may not track the target model’s behavior at scale (Li et al., 2025; Yu et al., 2024; Li et al., 2024). Other selection methods, including heuristics and LLM-based scoring, evaluate samples independently and ignore the underlying manifold structure (Zhuang et al., 2025; Xie et al., 2023a). This decoupling leads to structural blind spots. We either optimize macro-distribution without considering micro-quality, or filter samples at the cost of disrupting the hierarchical dependencies. This matters most for the **code corpus**, a rigid logical manifold defined by brittle syntax and hierarchical dependencies (Li et al., 2023; Guo et al., 2020; Feng, 2020). We thus ask: *How can we unify macro-distribution balancing and micro-quality selection to identify the "golden subset" within a structured code corpus?*

To bridge this gap, we introduce **UniGeM (Geometric Exploration and Mining)**, a framework that unifies macro-distribution balancing and micro-quality selection. Unlike existing methods that rely on external reference datasets for alignment or re-

quire training expensive proxy models to estimate data importance, UniGeM views these tasks as a unified *manifold approximation* problem. First, **Macro-Exploration (Stage-I)** discovers semantic manifolds via stability-driven clustering to optimize macro-distribution balancing. Subsequently, **Micro-Mining (Stage-II)** performs micro-quality selection using intrinsic geometric priors to capture structural and logical dependencies. By measuring manifold deviations, UniGeM distills a **golden subset** that serves as a faithful approximation of the data manifold, preserving global topology and local dependency structure. We validate UniGeM by training 8B and 16B Mixture-of-Experts (MoE) models (both with 1.4B active parameters) from scratch on a 100B-token code-and-text mixture. Results show that UniGeM achieves superior data efficiency and model performance compared to existing baselines. Our main contributions are:

- **Unified Framework:** We propose a geometry-centric framework to unify macro-distribution balancing and micro-quality selection via manifold approximation.
- **Proven Data Efficiency:** UniGeM achieves  $2.0\times$  **data efficiency** (vs. a random baseline) and performs better after one epoch.
- **Broader Gains:** UniGeM improves overall performance over SOTA baselines and shows stronger multilingual generalization.
- **Topological Insights:** Our ablations suggest that coverage and local structure matter for reasoning beyond per-sample quality scores, consistent with a manifold-approximation view of curation.

## 2 Proposed Method: The UniGeM Framework

In this section, we detail the implementation of UniGeM and provide a theoretical analysis to demonstrate effectiveness in data curation.

### 2.1 Problem Formulation and Overview

Let  $\mathcal{D}_{raw} = \{x_i\}_{i=1}^N$  denote the uncurated corpus. We employ an embedding model  $f_\theta$  to map each sample  $x_i$  to a normalized feature vector  $\mathbf{e}_i \in \mathbb{S}^{d-1}$ . We use this normalized space as a practical latent manifold  $\mathcal{M}$  for clustering and selection. We cast data selection not as distributional alignment to an external reference, but as a self-contained **manifold approximation** problem. Our objective is to identify a subset  $S \subset \mathcal{D}_{raw}$  that maximally preserves the topological structure of  $\mathcal{M}$  while suppressing

off-manifold outliers. Table 1 summarizes the key notations utilized throughout this framework.

Table 1: Summary of Key Notations and Symbols.

Symbol	Definition
<b>I. Macro-Exploration (Stage-I)</b>	
$C_k$	Global semantic cluster.
$\mathbf{z}_k, s_k$	Geometric features and score.
$\mathbf{w}$	Spectral consensus weights.
$K^*$	Optimal manifold resolution.
$T_{scale}$	Alignment scaling factor.
$r_k$	Global mixing budget.
<b>II. Micro-Mining (Stage-II)</b>	
$S_j$	Granular sub-cluster within $C_k$ .
$P_{S_j}$	Semantic scores via probes.
$\mathcal{L}_{struct}$	Structural penalty.
$\beta_{S_j}$	Geometric cohesion gate.
$W(S_j)$	Final sampling weight.
<b>III. Manifold Theory &amp; Approximation</b>	
$\mathcal{M}, d$	Latent manifold and dimensionality.
$\mathcal{E}(S)$	Wasserstein-2 approximation error.
$\sigma_k^2$	Intra-cluster variance of $C_k$ in latent space.
$\alpha_k$	Cluster mass (probability) $\mu(C_k)$ .
$\Delta_{gain}^{(k)}$	Variance mass rejected by Stage-II pruning.

### 2.2 Stage-I: Macro Exploration (Global Clustering and Weighting)

As illustrated in **Figure 2 (Left)**, Stage-I performs the **Macro-Distribution Balancing**. This phase transforms the raw manifold into structured semantic clusters and assigns a sampling budget to each cluster. The process consists of three sequential steps:

**Geometric Metrics and Scoring.** First, we characterize each candidate cluster  $C_k$  via a feature vector  $\mathbf{z}_k \in \mathbb{R}^4$ . These dimensions act as geometric proxies:

- **Cohesion ( $z_{coh}$ ):** Inverse intra-cluster distance. It prioritizes tight semantic structures akin to Neural Collapse (Papayan et al., 2020) to ensure gradient stability.
- **Cluster Size ( $z_{size}$ ):** Sample volume. Used to mitigate head redundancy (e.g., boilerplate) and shift focus to the information-dense long tail (Huang et al., 2024).
- **Sequence Length ( $z_{len}$ ):** The average token count, a proxy for verbosity, used to downweight overly long, low-information clusters.
- **Entropy ( $z_{ent}$ ):** Distributional impurity of language identifiers. It penalizes semantic ambiguity to ensure domain purity.

To unify these proxies, we combine the normalized signals with a weighted linear score to derive

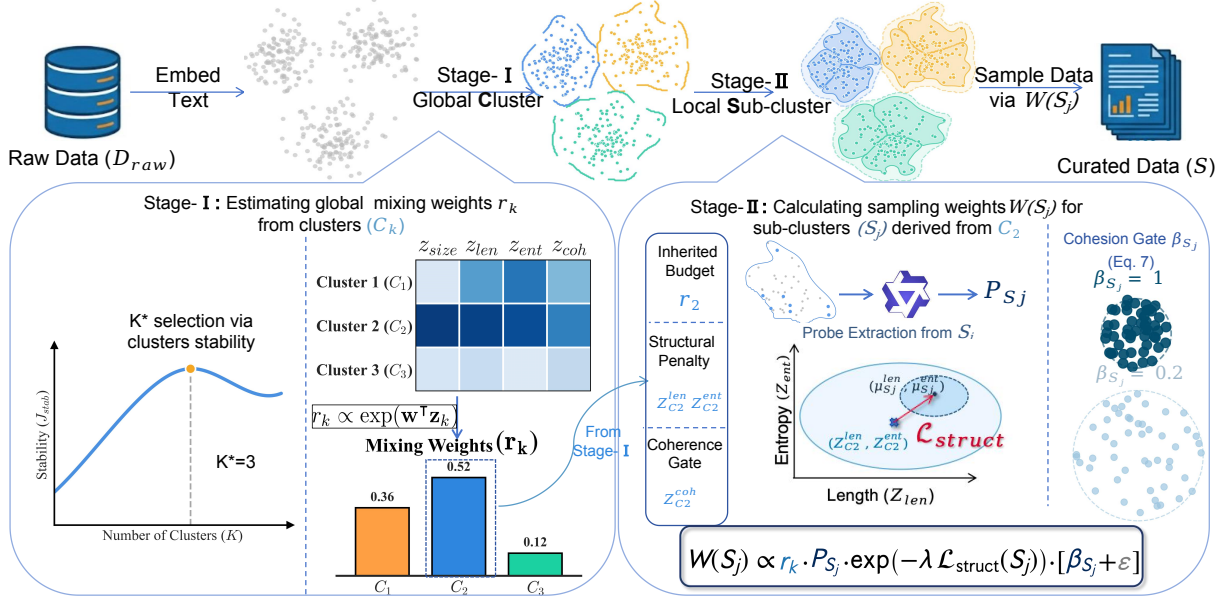


Figure 2: **Overview of the UniGeM framework.** The curation pipeline operates in two hierarchical stages: (Left) **Stage-I (Macro-Exploration)** identifies the intrinsic manifold resolution  $K^*$  via topological stability (Algorithm 1) and derives global mixing weights  $\mathbf{r}$  through a softmax over geometric scores (Eq. 5), ensuring comprehensive global coverage across diverse semantic regions. (Right) **Stage-II (Micro-Mining)** computes sub-cluster sampling weights  $W(S_j)$  by combining the inherited budget  $r_k$ , a semantic score modulated by a structural penalty  $\exp(-\lambda \mathcal{L}_{struct}(S_j))$  (Eq. 6), and a cohesion gate  $[\beta_{S_j} + \epsilon]$  (Eq. 7). This stage refines selection within each region to downweight off-manifold outliers while preserving representative local structural dependencies.

a scalar **Geometric Score**  $s_k$ . This score acts as a global quality metric, formulated as a contrast between structural coherence and statistical instability:

$$s_k = w_{\text{coh}} \tilde{z}_{\text{coh}} - \sum_{f \in \mathcal{F}_{\text{neg}}} w_f \tilde{z}_f \quad (1)$$

where  $\mathcal{F}_{\text{neg}} = \{\text{len}, \text{ent}, \text{size}\}$  denotes the set of negative factors, and  $\tilde{z}$  represents the Z-normalized magnitudes. The weighting vector  $\mathbf{w}$  is derived via Spectral Consensus (see Appendix A).

**Intrinsic Resolution Selection.** Central to this is a **Cross-Resolution Soft Alignment**. We construct a probabilistic bridge  $\pi_{K \rightarrow K'}$  based on centroid similarity:

$$\pi_{K \rightarrow K'}(j|i) = \text{Softmax} \left( T_{\text{scale}} \cdot \cos(\mu_i^{(K)}, \mu_j^{(K')}) \right) \quad (2)$$

Using this bridge, we project the geometric scores  $\mathbf{s}^{(K')}$  (derived in Eq. 1) from the finer resolution back to the current layer to obtain the **reconstructed scores**  $\hat{\mathbf{s}}^{(K)}$ :

$$\hat{s}_i^{(K)} = \sum_j \pi(j|i) s_j^{(K')} \quad (3)$$

### Algorithm 1 Intrinsic Resolution via Soft-Alignment

**Require:** Manifold  $\mathcal{M}$ , Range  $\mathcal{K}$ , Stride  $\Delta K$ .

**Ensure:** Optimal resolution  $K^*$ .

- 1:  $\mathcal{S} \leftarrow \emptyset$
- 2: **for**  $K \in \mathcal{K}$  **do**
- 3:  $K' \leftarrow K + \Delta K$
- 4: **1. Extract:** Compute centroids  $\mathcal{C}^{(K)}$  and scores  $\mathbf{s}^{(K)}$ .
- 5: **2. Align:** Construct soft bridge  $\pi_{K \rightarrow K'}$  (Eq. 2).
- 6: **3. Evaluate:** Calculate rank stability  $J_{\text{stab}}(K)$  (Eq. 4).
- 7:  $\mathcal{S} \leftarrow \mathcal{S} \cup \{(K, J_{\text{stab}}(K))\}$
- 8: **end for**
- 9: **return**  $K^* = \arg \max_K \{J : (K, J) \in \mathcal{S}\}$

Stability is quantified via rank correlation using a Kendall's  $\tau$  (detailed formulation in Appendix B):

$$J_{\text{stab}}(K) = \tau \left( \mathbf{s}^{(K)}, \hat{\mathbf{s}}^{(K)} \right) \quad (4)$$

The resolution  $K^*$  that maximizes this stability is selected as the intrinsic resolution (macro-granularity) via Algorithm 1.

**Sampling Budget Allocation.** With the optimal resolution  $K^*$  established, we finally allocate the

sampling budget  $\mathbf{r}$ . Leveraging the geometric scores from Eq. (1), the probability for each cluster is computed via a standard softmax:

$$r_k = \frac{\exp(s_k)}{\sum_{j=1}^{K^*} \exp(s_j)} \quad (5)$$

This framework ensures the budget concentrates on semantically dense regions (high  $s_k$ ).

### 2.3 Stage-II: Micro-Mining (Local Sub-Clustering and Refinement)

As illustrated in Figure 2 (Right), Stage-II performs the **Micro-Quality Selection** via local sub-clustering. This phase decomposes each global cluster  $C_k$  into fine-grained sub-clusters  $\{S_j\}$  to capture semantic diversity and topological details. The selection is refined through three coupled mechanisms:

**Probe-based Semantic Scoring.** To assess semantic content efficiently, we extract a small representative **Probe Set** (indicated by the solid centroids in Figure 2) from each sub-cluster  $S_j$ . We utilize an LLM as a **Knowledge Probe** to derive a **Semantic Score**  $P_{S_j}$ . Implementation details are provided in Appendix D.

**Relative Structural Consistency.** To enforce manifold consistency, we define a **Structural Penalty**  $\mathcal{L}_{struct}$  using a rectified Mahalanobis distance. It penalizes sub-clusters  $S_j$  where the empirical mean  $z_f(S_j)$  of features (e.g., length, entropy) exceeds the geometric consensus  $(\mu_f, \sigma_f)$  of the parent cluster  $C_k$ :

$$\mathcal{L}_{struct} = \sum_{f \in \{len, ent\}} \left[ \frac{z_f(S_j) - \mu_f^{(C_k)}}{\sigma_f^{(C_k)}} \right]_+^2 \quad (6)$$

We introduce  $\mathcal{L}_{struct}$  not only to ensure that the curated subset preserves the underlying manifold topology but also to mitigate the score saturation observed in semantic scores (Fig. 10 and Fig. 11).

**Geometric Cohesion Gate.** To ensure topological robustness, we apply a **Geometric Cohesion Gate**  $\beta_{S_j}$ , depicted as the gating module in Figure 2. Leveraging the cohesion metric  $z_{coh}$  defined in Stage-I, we modulate sampling confidence based on the sub-cluster’s compactness relative to its parent cluster:

$$\beta_{S_j} = \text{Sigmoid} \left( z_{coh}^{(S_j)} - z_{coh}^{(C_k)} \right) \quad (7)$$

This gating mechanism suppresses sub-clusters with lower cohesion than their parent ( $z_{coh}^{(S_j)} < z_{coh}^{(C_k)}$ ) while retaining structurally compact regions.

**Hierarchical Sampling Weight.** The final sampling probability  $W(S_j)$  synthesizes the global budget  $r_{C_k}$  with local metrics via a multiplicative modulation. It is computed as:

$$W(S_j) \propto \underbrace{r_k}_{\text{Budget}} \cdot \underbrace{P_{S_j} \exp(-\lambda \mathcal{L}_{struct}(S_j))}_{\text{Geometry-Aware Score}} \cdot \underbrace{[\beta_{S_j} + \epsilon]}_{\text{Gate}} \quad (8)$$

Within the geometry-aware score, we use a semantic score modulated by a geometric penalty to suppress off-manifold sub-clusters. Simultaneously, the gate  $[\beta_{S_j} + \epsilon]$  ensures cohesion, with  $\epsilon$  preventing mode collapse by maintaining a minimal exploration floor.

### 2.4 Theoretical Analysis: Manifold Approximation

We frame data selection as minimizing the transport cost between the empirical measure  $\hat{\mu}_S$  and the true manifold distribution  $\mu$  (Villani, 2009). Let  $f_\theta : \mathcal{X} \rightarrow \mathbb{R}^d$  be an  $L$ -Lipschitz embedding into a latent manifold of intrinsic dimension  $d \ll D$  (Pope et al., 2021; Du et al., 2021).

**Definition 1 (Approximation Error).** We quantify the quality of subset  $S$  by the squared Type-2 Wasserstein distance (Peyré et al., 2019):

$$\begin{aligned} \mathcal{E}(S) &\triangleq W_2^2(\mu, \hat{\mu}_S) \\ &= \inf_{\gamma \in \Pi(\mu, \hat{\mu}_S)} \iint_{\mathcal{M}^2} \|x - y\|^2 d\gamma(x, y). \end{aligned} \quad (9)$$

A constructive two-stage transport argument (Appendix C) yields the following constant-factor decomposition:

$$\begin{aligned} \mathcal{E}(S) &\leq 2 \underbrace{\sum_{k=1}^K \int_{C_k} \|x - \mathbf{c}_k\|^2 d\mu(x)}_{\text{Stage-I: Quantization}} \\ &\quad + 2 \underbrace{\sum_{k=1}^K \alpha_k \mathbb{E}_{x \sim \hat{\mu}_{S_k}} \|x - \mathbf{c}_k\|^2}_{\text{Stage-II: Pruning}}. \end{aligned} \quad (10)$$

**Theorem 1 (UniGeM Bound).** For density  $p(x)$ , the UniGeM approximation error satisfies:

$$\mathcal{E}(S_{UniGeM}) \leq 2C_d K^{-2/d} + 2 \sum_{k=1}^K \alpha_k (\sigma_k^2 - \Delta_{gain}^{(k)}), \quad (11)$$

where  $C_d$  is the Zador constant (Zador, 1982),  $\alpha_k \triangleq \mu(C_k)$ , and  $\sigma_k^2 \triangleq \mathbb{E}_{x \sim \mu(\cdot|C_k)} \|x - c_k\|^2$ . Here  $\Delta_{gain}^{(k)} \geq 0$  captures the within-cluster reduction induced by Stage-II pruning.  $\mathcal{L}_{struct}$  is a practical signal for identifying the outliers driving  $\Delta_{gain}^{(k)}$ . **Remark.** The bound is derived under idealized modeling assumptions and is intended to guide the design rather than exactly characterize all engineering details.

### 3 Experimental Setup

We evaluate UniGeM on **large-scale code pre-training**. Code is a demanding setting for data curation: small syntax issues can break execution (Li et al., 2023), the corpus spans many programming languages (Feng, 2020), and programs follow hierarchical dependencies (Guo et al., 2020). These properties make naive filtering brittle, and they let us test whether UniGeM preserves both coverage and local structure.

#### 3.1 Corpus Construction and Sampling Protocol

We construct a 100B-token training corpus by mixing **The Stack Dedup** (Kocetkov et al., 2022) and Common Crawl (Schäfer, 2017) in a fixed 7:3 code-to-text ratio. The 30B text component is frozen across all experiments. This fixed ratio ensures that gains on code benchmarks stem from UniGeM’s strategic data blending rather than a simple inflation of total code volume (Xie et al., 2023b).

#### 3.2 UniGeM Implementation Details

Deploying geometric clustering on terabytes of data presents a computational challenge. We adopt a *Probing-and-Scaling* strategy:

- **Phase I: Macro-Topology Discovery.** We first removed near-duplicates of evaluation benchmarks from the raw corpus. Subsequently, we embed the corpus using *Qwen3-embedding* (Zhang et al., 2025b). To ensure scalability, we perform K-means clustering on a representative subset ( $\mathcal{D}_{probe} \approx 20\%$ ) to identify the optimal macro-granularity  $K^* = 72$ . Each sample in the full corpus is assigned to its nearest global centroid.
- **Phase II: Hierarchical Geometric Mining.** Within each global cluster  $C_k$ , we execute localized sub-clustering to capture fine-grained micro-structures. The sub-cluster density  $S_j$

is determined by the square root of the cluster population.

#### 3.3 Model and Hyperparameter

We adopt a fine-grained sparse Mixture-of-Experts (MoE) Transformer, increasing the total expert capacity while keeping the number of activated parameters per token fixed to preserve inference efficiency (Liu et al., 2024). We instantiate **UniGeM-8B** (32 experts) and **UniGeM-16B** (64 experts), both with **1.4B** active parameters. All models are trained **from scratch** with the same training recipe; only the data curation differs<sup>1</sup>. For geometric selection, we use  $\lambda = 0.5$ ,  $T_{scale} = 20$ , and  $\epsilon = 0.01$ . See Appendix E for details.

#### 3.4 Compared Methods

##### Baseline & SOTA Methods

- **Random Sampling:** Uniform sampling from the raw corpus as a reference baseline.
- **Meta-rater** (Zhuang et al., 2025): A representative *instance-level selection* method relying on *LLM-based scoring* to rank individual data quality.
- **Nemotron-CLIMB** (Diao et al., 2025): A representative *domain-level data mixing* method that optimizes global balancing weights.

To ensure competitive baselines, we apply a **code adaptation** to Meta-rater and CLIMB to reduce domain mismatch (Appendix E).

**Ablation Variants** To validate the hierarchical design, we compare UniGeM against:

- **Cluster Random:** Replaces the coarse-to-fine hierarchy with a single flat clustering stage ( $K = K_{total}$ ) and performs random sampling with a uniform ratio in each cluster.
- **w/o Hierarchy:** Uses the same flat clustering structure as above ( $K = K_{total}$ ) but selects data via UniGeM’s micro-mining metrics instead of random sampling.
- **w/o Stage-I:** Assigns uniform global budgets ( $r_k = 1/K$ ) while retaining Stage-II mining, utilizing **global feature statistics** as the reference for both the structural and the geometric cohesion gate.
- **w/o Stage-II:** Uses optimized global weights  $\mathbf{r}$  but performs random sampling within clusters.

<sup>1</sup>Reproducibility statement is discussed in Appendix H.

Table 2: Coding benchmark results. We compare the scaling properties between UniGeM-16B and UniGeM-8B under a constant inference budget (1.4B active parameters). The best scores in Block 2 and Block 3 are underlined.

Setting		Avg.	HE	HE <sup>+</sup>	MBPP	MBPP <sup>+</sup>	LiveCode	CruxEval	
Methods	Epoch	Score	Pass@1	Pass@1	Pass@1	Pass@1	Pass@1	Input	Output
<b>UniGeM-16B (16B Total / 1.4B Active)</b>									
Random	1.0	32.9	55.9	48.6	33.5	37.6	1.1	29.8	24.1
<b>UniGeM (Ours)</b>	0.5	32.3 <span style="color: red;">↓0.6</span>	49.9 <span style="color: red;">↓6.0</span>	43.8 <span style="color: red;">↓4.8</span>	34.6 <span style="color: green;">↑1.1</span>	40.9 <span style="color: green;">↑3.3</span>	6.4 <span style="color: green;">↑5.3</span>	26.0 <span style="color: red;">↓3.8</span>	24.6 <span style="color: green;">↑0.5</span>
<b>UniGeM (Ours)</b>	1.0	39.5 <span style="color: green;">↑6.6</span>	59.3 <span style="color: green;">↑3.4</span>	55.1 <span style="color: green;">↑6.5</span>	41.8 <span style="color: green;">↑8.3</span>	47.8 <span style="color: green;">↑10.2</span>	10.8 <span style="color: green;">↑9.7</span>	31.6 <span style="color: green;">↑1.8</span>	30.1 <span style="color: green;">↑6.0</span>
<b>UniGeM-8B (8B Total / 1.4B Active)</b>									
Random	1.0	29.1	45.3	44.1	30.5	34.2	2.0	22.3	25.5
Meta-rater	1.0	35.0 <span style="color: green;">↑5.9</span>	53.2 <span style="color: green;">↑7.9</span>	49.6 <span style="color: green;">↑5.5</span>	35.3 <span style="color: green;">↑4.8</span>	43.5 <span style="color: green;">↑9.3</span>	4.7 <span style="color: green;">↑2.7</span>	<u>31.4</u> <span style="color: green;">↑9.1</span>	27.7 <span style="color: green;">↑2.2</span>
CLIMB	1.0	35.2 <span style="color: green;">↑6.1</span>	52.4 <span style="color: green;">↑7.1</span>	48.8 <span style="color: green;">↑4.7</span>	36.7 <span style="color: green;">↑6.2</span>	45.5 <span style="color: green;">↑11.3</span>	6.9 <span style="color: green;">↑4.9</span>	29.7 <span style="color: green;">↑7.4</span>	26.6 <span style="color: green;">↑1.1</span>
<b>UniGeM (Ours)</b>	0.5	30.0 <span style="color: green;">↑0.9</span>	47.7 <span style="color: green;">↑2.4</span>	43.5 <span style="color: red;">↓0.6</span>	32.4 <span style="color: green;">↑1.9</span>	42.2 <span style="color: green;">↑8.0</span>	5.2 <span style="color: green;">↑3.2</span>	16.6 <span style="color: red;">↓5.7</span>	22.5 <span style="color: red;">↓3.0</span>
<b>UniGeM (Ours)</b>	1.0	36.4 <span style="color: green;">↑7.3</span>	53.7 <span style="color: green;">↑8.4</span>	50.3 <span style="color: green;">↑6.2</span>	37.4 <span style="color: green;">↑6.9</span>	46.6 <span style="color: green;">↑12.4</span>	7.8 <span style="color: green;">↑5.8</span>	31.0 <span style="color: green;">↑8.7</span>	28.4 <span style="color: green;">↑2.9</span>
<b>Ablation Study (based on UniGeM-8B)</b>									
Random	1.0	29.1	45.3	44.1	30.5	34.2	2.0	22.3	25.5
Cluster Random	1.0	29.2 <span style="color: green;">↑0.1</span>	46.0 <span style="color: green;">↑0.7</span>	44.1	31.1 <span style="color: green;">↑0.6</span>	34.6 <span style="color: green;">↑0.4</span>	1.8 <span style="color: red;">↓0.2</span>	21.4 <span style="color: red;">↓0.9</span>	25.2 <span style="color: red;">↓0.3</span>
w/o Hierarchy	1.0	33.0 <span style="color: green;">↑3.9</span>	51.4 <span style="color: green;">↑6.1</span>	48.1 <span style="color: green;">↑4.0</span>	34.5 <span style="color: green;">↑4.0</span>	40.1 <span style="color: green;">↑5.9</span>	1.9 <span style="color: red;">↓0.1</span>	27.0 <span style="color: green;">↑4.7</span>	28.0 <span style="color: green;">↑2.5</span>
w/o Stage-II	1.0	31.2 <span style="color: green;">↑2.1</span>	45.8 <span style="color: green;">↑0.5</span>	43.8 <span style="color: red;">↓0.3</span>	33.8 <span style="color: green;">↑3.3</span>	40.1 <span style="color: green;">↑5.9</span>	2.3 <span style="color: green;">↑0.3</span>	24.3 <span style="color: green;">↑2.0</span>	28.3 <span style="color: green;">↑2.8</span>
w/o Stage-I	1.0	34.9 <span style="color: green;">↑5.8</span>	52.2 <span style="color: green;">↑6.9</span>	49.5 <span style="color: green;">↑5.4</span>	36.1 <span style="color: green;">↑5.6</span>	43.1 <span style="color: green;">↑8.9</span>	4.5 <span style="color: green;">↑2.5</span>	<u>31.3</u> <span style="color: green;">↑9.0</span>	27.5 <span style="color: green;">↑2.0</span>
<b>UniGeM (Ours)</b>	1.0	36.4 <span style="color: green;">↑7.3</span>	53.7 <span style="color: green;">↑8.4</span>	50.3 <span style="color: green;">↑6.2</span>	37.4 <span style="color: green;">↑6.9</span>	46.6 <span style="color: green;">↑12.4</span>	7.8 <span style="color: green;">↑5.8</span>	31.0 <span style="color: green;">↑8.7</span>	28.4 <span style="color: green;">↑2.9</span>

### 3.5 Code Benchmarks

We evaluate UniGeM on a diverse set of coding benchmarks covering code generation, robustness/reasoning, and multilingual proficiency:

- **Code Generation:** We utilize *HumanEval(+)* (HE) (Chen, 2021) and *MBPP(+)* (Austin et al., 2021) to evaluate generation performance.
- **Robustness & Reasoning:** We assess temporal generalization using *LiveCodeBench (LCB)* (Jain et al., 2024) and execution prediction capabilities using *CruxEval* (Gu et al., 2024).
- **Multilingual Proficiency:** We conduct a fine-grained analysis across multiple programming languages using *MultiPL-E* (Cassano et al., 2023).

## 4 Result Analysis

In this section, we evaluate UniGeM from multiple perspectives. We compare it with state-of-the-art baselines, examine cross-lingual generalization, and study the impact of the hierarchical design and training efficiency through ablation and scaling

analyses.

### 4.1 Comparison with State-of-the-Art

We compare UniGeM against the Random baseline, **CLIMB** and **Meta-rater** on the 8B MoE model. Table 2 summarizes the results; the **Avg.** score represents the unweighted arithmetic mean of all reported metrics.

#### Overall Superiority and Data Efficiency.

UniGeM achieves the best overall score (**36.4**), above the adapted domain-mixing baseline (**CLIMB**, 35.2) and the adapted instance-selection baseline (**Meta-rater**, 35.0). It is also more data-efficient: UniGeM reaches **30.0** after 0.5 epochs, surpassing Random sampling at 1.0 epoch (29.1), corresponding to an approximate 2.0× efficiency gain. After 1.0 epoch, UniGeM further improves to **36.4**, indicating that the gain is not limited to early training.

**Reasoning and Generalization.** UniGeM’s advantage is particularly robust on complex execution and out-of-distribution (OOD) tasks. On

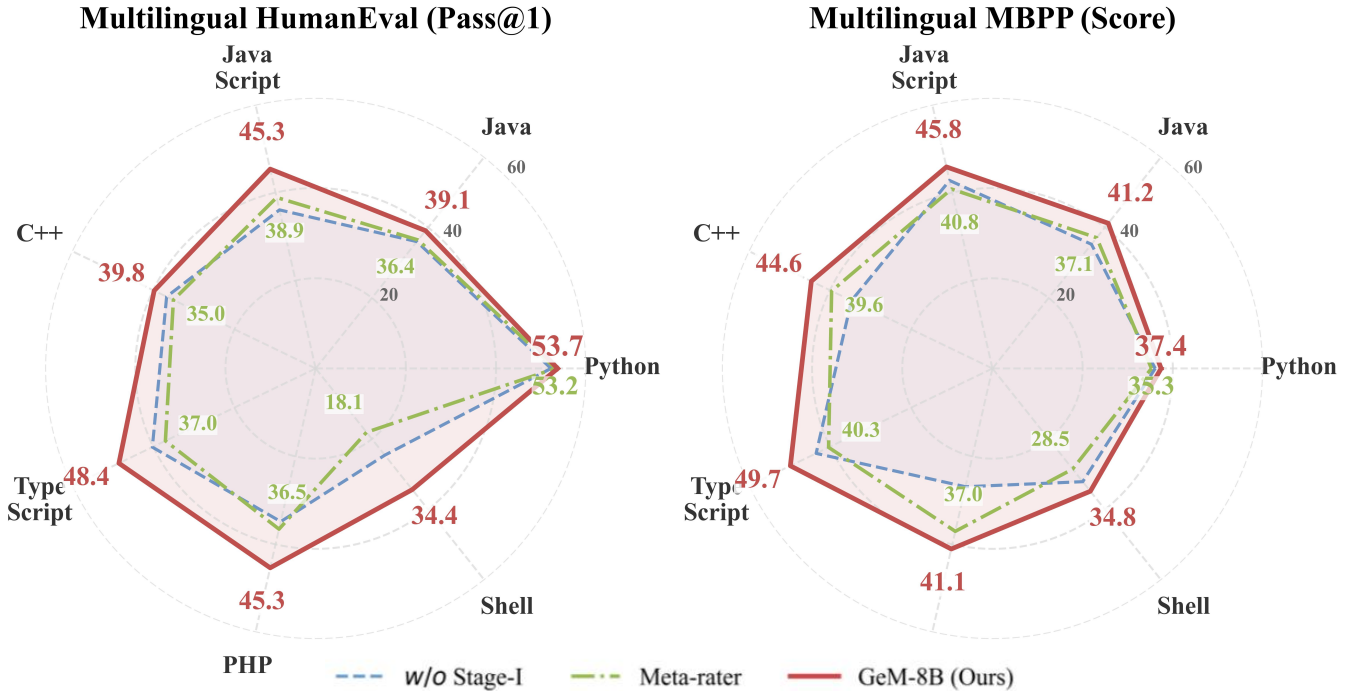


Figure 3: **Impact of Macro-Exploration (Stage-I).** UniGeM (red) outperforms Meta-rater (green) and the *w/o Stage-I* variant (blue) across 7 languages, indicating that global clustering improves multilingual generalization.

**LiveCodeBench**, UniGeM scores **7.8**, clearly surpassing CLIMB (6.9) and significantly exceeding Meta-rater (4.7). On **CruxEval**, the adapted Meta-rater remains highly competitive and slightly leads on **Input** prediction, suggesting that code-aware LLM scoring captures common logical patterns. UniGeM performs best on **Output** execution prediction (**CruxEval-Output**) with **28.4**, indicating better coverage of more complex execution behaviors.

## 4.2 Multilingual Proficiency

We evaluate cross-lingual generalization on **Multilingual HumanEval** and **MBPP** (Fig. 3). UniGeM consistently dominates the Meta-rater and *w/o Stage-I* variant across 7 languages. On HumanEval, UniGeM achieves decisive leads in low-resource domains like **Shell** (34.4% vs. Meta’s 18.1%) and **TypeScript** (48.4% vs. 37.0%), while maintaining clear advantages in strict-syntax languages like **C++** (39.8% vs. 35.0%) and **Java** (39.1% vs. 36.4%). This superiority persists on MBPP (e.g., 44.6% vs. 39.6% in C++), confirming UniGeM’s robust, language-agnostic transfer capabilities.

## 4.3 Ablation and Hyperparameter Analysis

We conduct ablations and sensitivity analyses on UniGeM-8B to isolate the contribution of each component.

**Contribution of Model Components.** Ablations highlight the role of each stage. Removing Stage-I (*w/o Stage-I*) replaces learned global budgets with uniform allocation, which weakens global coverage; performance stays relatively stable on high-resource languages (e.g., Python) but drops on under-represented languages such as Shell and PHP in multilingual evaluation. Conversely, removing Stage-II (*w/o Stage-II*) reduces HumanEval Pass@1 from 53.7% to 45.8%, indicating that local geometric refinement is important even with optimized macro mixing.

**Statistical Properties of Stage-I Features.** Statistical analysis confirms that extensive features (e.g., length) follow a Log-Normal distribution while intensive features (e.g., cohesion) are naturally stable, supporting our hybrid geometric priors (details in Appendix F).

**Sensitivity to Global Granularity ( $K$ ).** The cluster number  $K$  controls the resolution of the global approximation. As shown in Figure 4, the stability index  $J_{stab}(K)$  increases and then plateaus beyond  $K \approx 60$ . We choose  $K^* = 72$  as a stable setting with finer semantic resolution; smaller  $K$  (e.g.,  $K < 40$ ) suffers from under-segmentation, merging distinct semantic domains and reducing stability. Beyond  $K^*$ , the marginal gain ( $\Delta J$ ) diminishes to near zero, suggesting that

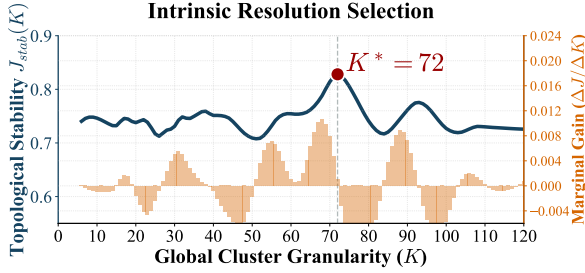


Figure 4: **Intrinsic Resolution Selection.** The stability index  $J_{\text{stab}}(K)$  (blue line) rises and then plateaus, attaining its maximum at  $K^* = 72$ . The marginal gain (orange bars) diminishes significantly beyond this point.

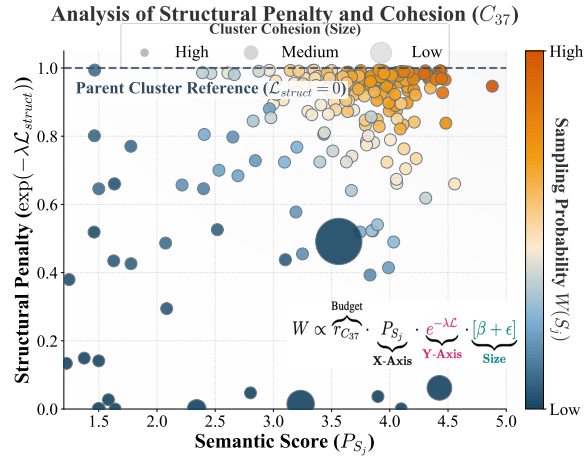


Figure 5: Empirical distribution of a representative subset of sub-clusters in  $C_{37}$ . Sub-clusters are mapped by **Semantic Score** ( $P_{S_j}$ , X-axis) and **Structural Penalty** (Y-axis), with bubble size representing **Cluster Cohesion**. The color gradient indicates the **Sampling Probability**  $W(S_j)$ , illustrating the budget allocation toward sub-clusters.

72 clusters are sufficient to capture the corpus’s semantic structure without over-segmentation or fragmentation.

#### Analysis of Structural Penalty and Cohesion.

Figure 5 illustrates how UniGeM balances semantic utility and structural consistency at the sub-cluster level. Sub-clusters with higher **Semantic Score** and a small structural deviation (i.e., **Structural Penalty** close to 1) receive higher sampling probability, whereas sub-clusters that score well semantically but deviate structurally are downweighted. **Cluster Cohesion** (bubble size) further biases sampling toward compact, well-formed sub-clusters that better represent the parent cluster.

#### 4.4 Scaling Analysis

**Model Scaling Dynamics.** We further study scaling by training a larger **UniGeM-16B** model. At 0.5 epochs, UniGeM-16B reaches **32.3**, which is lower than its 1.0-epoch score, consistent with

larger models being more data-hungry (Hoffmann et al., 2022). After 1.0 epoch, UniGeM-16B improves to **39.5 (+6.6 over Random)**, showing that the curated corpus continues to benefit training as model size increases.

## 5 Related Work

Current work on data efficiency is moving beyond static filtering toward methods that model dataset structure. One prominent line trains **proxy models** to derive mixing or importance weights (e.g., DoReMi (Xie et al., 2023b) and DCLM (Li et al., 2024)), but this can be compute-heavy and may introduce “proxy bias,” where signals from small proxies do not transfer cleanly to larger target models (Mindermann et al., 2022; Sorscher et al., 2022). A related family relies on reference datasets for alignment or selection (Xie et al., 2023a; Li et al., 2025), while model-aware approaches such as Mates (Yu et al., 2024; Zhang et al., 2025a) use influence-style estimates to capture sample-level contributions. To better preserve reasoning ability, QuaDMix (Liu et al., 2025) explicitly balances quality and diversity, though gradient- or training-intensive signals can limit scalability (Mindermann et al., 2022; Li et al., 2024). More principled directions aim to avoid expensive training signals altogether: DDOQ (Tan and Slade, 2025) casts selection as pushforward optimal quantization, improving over heuristic clustering schemes (Chen et al., 2023; Diao et al., 2025), and Wasserstein-manifold views model dataset dynamics in a way that goes beyond flat domain mixing toward preserving structure relevant for complex reasoning (Atanackovic et al., 2024).

## 6 Conclusion

We introduced **UniGeM**, a hierarchical framework that unifies macro-distribution balancing and micro-quality selection through *manifold approximation*. By using topological stability to choose the global resolution and geometric priors for instance mining, UniGeM curates a compact, structure-preserving training set from code corpora. Experiments with 8B and 16B MoE models show **2.0× data efficiency** over a random baseline and **better one-epoch performance** than strong adapted baselines, with consistent gains in code reasoning and multi-lingual evaluations.

## 504 Limitations

505 Despite its effectiveness, this work has several lim-  
506 itations:

- 507 1. **Domain Specificity:** Our evaluation focused  
508 primarily on the **code corpus**. While code pro-  
509 vides a rigorous testbed for geometric structures,  
510 the efficacy of UniGeM on massive, heteroge-  
511 neous **general web text** mixtures remains to be  
512 fully explored.
- 513 2. **Computational Overhead:** The initial global  
514 embedding and clustering phase, while miti-  
515 gated by our probing-and-scaling strategy, still  
516 requires non-trivial resources when applied to  
517 trillion-token scales.
- 518 3. **Static Pipeline:** The current framework oper-  
519 ates as a pre-processing step. Future work is re-  
520 quired to integrate UniGeM into online training  
521 pipelines to allow for dynamic, cross-domain  
522 manifold updates as the model’s data needs  
523 evolve.

## 524 References

525 Lazar Atanackovic, Xi Zhang, Brandon Amos, Mathieu  
526 Blanchette, Leo J Lee, Yoshua Bengio, Alexander  
527 Tong, and Kirill Neklyudov. 2024. Meta flow match-  
528 ing: Integrating vector fields on the wasserstein man-  
529 ifold. *arXiv preprint arXiv:2408.14608*.

530 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten  
531 Bosma, Henryk Michalewski, David Dohan, Ellen  
532 Jiang, Carrie Cai, Michael Terry, Quoc Le, and 1  
533 others. 2021. Program synthesis with large language  
534 models. *arXiv preprint arXiv:2108.07732*.

535 Federico Cassano, John Gouwar, Daniel Nguyen, Syd-  
536 ney Nguyen, Luna Phipps-Costin, Donald Pinckney,  
537 Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson,  
538 Molly Q Feldman, and 1 others. 2023. Multipl-e: A  
539 scalable and polyglot approach to benchmarking neu-  
540 ral code generation. *IEEE Transactions on Software  
541 Engineering*, 49(7):3675–3691.

542 Mark Chen. 2021. Evaluating large language models  
543 trained on code. *arXiv preprint arXiv:2107.03374*.

544 Mayee Chen, Nicholas Roberts, Kush Bhatia, Jue Wang,  
545 Ce Zhang, Frederic Sala, and Christopher Ré. 2023.  
546 Skill-it! a data-driven skills framework for under-  
547 standing and training language models. *Advances in  
548 Neural Information Processing Systems*, 36:36000–  
549 36040.

550 Shizhe Diao, Yu Yang, Yonggan Fu, Xin Dong, Dan  
551 Su, Markus Kliegl, Zijia Chen, Peter Belcak, Yoshi  
552 Suhara, Hongxu Yin, and 1 others. 2025. Climb:  
553 Clustering-based iterative data mixture bootstrap-  
554 ping for language model pre-training. *arXiv preprint  
555 arXiv:2504.13161*.

556 Yilun Du, Katie Collins, Josh Tenenbaum, and Vincent  
557 Sitzmann. 2021. Learning signal-agnostic manifolds  
558 of neural fields. *Advances in Neural Information  
559 Processing Systems*, 34:8320–8331.

560 Z Feng. 2020. Codebert: A pre-trained model for  
561 program-ming and natural languages. *arXiv preprint  
562 arXiv:2002.08155*.

563 Samir Yitzhak Gadre, Georgios Smyrnis, Vaishaal  
564 Shankar, Suchin Gururangan, Mitchell Wortsman,  
565 Rulin Shao, Jean Mercat, Alex Fang, Jeffrey Li,  
566 Sedrick Keh, and 1 others. 2024. Language models  
567 scale reliably with over-training and on downstream  
568 tasks. *arXiv preprint arXiv:2403.08540*.

569 Robert M. Gray and David L. Neuhoff. 2002. Quan-  
570 tization. *IEEE transactions on information theory*,  
571 44(6):2325–2383.

572 Alex Gu, Baptiste Rozière, Hugh Leather, Armando  
573 Solar-Lezama, Gabriel Synnaeve, and Sida I Wang.  
574 2024. **Cruxeval: A benchmark for code reason-  
575 ing, understanding and execution**. *arXiv preprint  
576 arXiv:2401.03065*.

577 Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio  
578 César Teodoro Mendes, Allie Del Giorno, Sivakanth  
579 Gopi, Mojan Javaheripi, Piero Kauffmann, Gus-  
580 tavo de Rosa, Olli Saarikivi, and 1 others. 2023.  
581 Textbooks are all you need. *arXiv preprint  
582 arXiv:2306.11644*.

583 Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu  
584 Tang, Shujie Liu, Long Zhou, Nan Duan, Alexey Svy-  
585 atkovskiy, Shengyu Fu, and 1 others. 2020. Graph-  
586 codebert: Pre-training code representations with data  
587 flow. *arXiv preprint arXiv:2009.08366*.

588 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch,  
589 Elena Buchatskaya, Trevor Cai, Eliza Rutherford,  
590 Diego de Las Casas, Lisa Anne Hendricks, Johannes  
591 Welbl, Aidan Clark, and 1 others. 2022. Training  
592 compute-optimal large language models. *Advances  
593 in Neural Information Processing Systems*, 35:30016–  
594 30030.

595 Jing Huang, Diyi Yang, and Christopher Potts. 2024.  
596 Demystifying verbatim memorization in large lan-  
597 guage models. *arXiv preprint arXiv:2407.17817*.

598 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia  
599 Yan, Tianjun Zhang, Sida Wang, Armando Solar-  
600 Lezama, Koushik Sen, and Ion Stoica. 2024. Live-  
601 codebench: Holistic and contamination free eval-  
602 uation of large language models for code. *arXiv  
603 preprint arXiv:2403.07974*.

604 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B  
605 Brown, Benjamin Chess, Rewon Child, Scott Gray,  
606 Alec Radford, Jeffrey Wu, and Dario Amodei. 2020.  
607 Scaling laws for neural language models. *arXiv  
608 preprint arXiv:2001.08361*.

609	Denis Kocetkov, Raymond Li, Loubna Ben Allal, Jia Li, Chenghao Mou, Carlos Muñoz Ferrandis, Yacine Jernite, Margaret Mitchell, Sean Hughes, Thomas Wolf, and 1 others. 2022. The stack: 3 tb of permissively licensed source code. <i>arXiv preprint arXiv:2211.15533</i> .	665
610		666
611		667
612		
613		
614		
615	Michel Ledoux. 2001. <i>The concentration of measure phenomenon</i> . 89. American Mathematical Soc.	
616		
617	Jiazheng Li, Lu Yu, Qing Cui, Zhiqiang Zhang, Jun Zhou, Yanfang Ye, and Chuxu Zhang. 2025. <a href="#">MASS: Mathematical data selection via skill graphs for pretraining large language models</a> . <i>arXiv preprint arXiv:2503.14917</i> .	
618		
619		
620		
621		
622	Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, and 1 others. 2023. Starcoder: may the source be with you! <i>arXiv preprint arXiv:2305.06161</i> .	
623		
624		
625		
626		
627	Ruihang Li, Yixuan Wei, Miaosen Zhang, Nenghai Yu, Han Hu, and Houwen Peng. 2024. <a href="#">Scalingfilter: Assessing data quality through inverse utilization of scaling laws</a> . <i>arXiv preprint arXiv:2408.08310</i> .	
628		
629		
630		
631	Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Deng, Chong Ruan, Damai Dai, Daya Guo, and 1 others. 2024. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model. <i>arXiv preprint arXiv:2405.04434</i> .	
632		
633		
634		
635		
636		
637	Fengze Liu, Weidong Zhou, Binbin Liu, Zhimiao Yu, Yifan Zhang, Haobin Lin, Yifeng Yu, Bingni Zhang, Xiaohuan Zhou, Taifeng Wang, and 1 others. 2025. Quadmix: Quality-diversity balanced data selection for efficient llm pretraining. <i>arXiv preprint arXiv:2504.16511</i> .	
638		
639		
640		
641		
642		
643	Sören Mindermann, Jan M Brauner, Muhammed Razzak, Mrinank Sharma, Andreas Kirsch, Winnie Xu, Benedikt Höltgen, Aidan N Gomez, Adrien Morisot, Sebastian Farquhar, and 1 others. 2022. Prioritized training on points that are learnable, worth learning, and not yet learnt (rho-loss). In <i>International Conference on Machine Learning</i> , pages 15630–15660. PMLR.	
644		
645		
646		
647		
648		
649		
650		
651	Vardan Papyan, XY Han, and David L Donoho. 2020. Prevalence of neural collapse during the terminal phase of deep learning training. <i>Proceedings of the National Academy of Sciences</i> , 117(40):24652–24663.	
652		
653		
654		
655		
656	Gabriel Peyré, Marco Cuturi, and 1 others. 2019. Computational optimal transport: With applications to data science. <i>Foundations and Trends® in Machine Learning</i> , 11(5-6):355–607.	
657		
658		
659		
660	Phil Pope, Chen Zhu, Ahmed Abdelkader, Micah Goldblum, and Tom Goldstein. 2021. The intrinsic dimension of images and its impact on learning. In <i>International Conference on Learning Representations (ICLR)</i> .	
661		
662		
663		
664		
	Roland Schäfer. 2017. Accurate and efficient general-purpose boilerplate detection for crawled web corpora. <i>Lang. Resour. Eval.</i> , 51(3):873–889.	665
		666
		667
	ByteDance Seed, Yuyu Zhang, Jing Su, Yifan Sun, Chenguang Xi, Xia Xiao, Shen Zheng, Anxiang Zhang, Kaibo Liu, Daoguang Zan, and 1 others. 2025. Seed-coder: Let the code model curate data for itself. <i>arXiv preprint arXiv:2506.03524</i> .	668
		669
		670
		671
		672
	Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Morcos Ariyo. 2022. Beyond neural scaling laws: beating power law scaling via data pruning. <i>Advances in Neural Information Processing Systems</i> , 35:19523–19536.	673
		674
		675
		676
		677
	Hong Ye Tan and Emma Slade. 2025. Dataset distillation as pushforward optimal quantization. <i>arXiv preprint arXiv:2501.07681</i> .	678
		679
		680
	Pablo Villalobos and 1 others. 2022. Will we run out of data? an analysis of the limits of scaling datasets in machine learning. <i>arXiv preprint arXiv:2211.04325</i> .	681
		682
		683
	Cédric Villani. 2009. <i>Optimal transport: old and new</i> , volume 338. Springer.	684
		685
	Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy Liang, Quoc V Le, Tengyu Ma, and Adams Wei Yu. 2023a. Data selection for language models via importance resampling. <i>arXiv preprint arXiv:2302.03169</i> .	686
		687
		688
		689
		690
	Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy S Liang, Quoc V Le, Tengyu Ma, and Adams Wei Yu. 2023b. Doremi: Optimizing data mixtures speeds up language model pretraining. In <i>NeurIPS</i> .	691
		692
		693
		694
		695
	Zichun Yu, Spandan Das, and Chenyan Xiong. 2024. Mates: Model-aware data selection for efficient pretraining with data influence models. <i>Advances in Neural Information Processing Systems</i> , 37:108735–108759.	696
		697
		698
		699
		700
	Paul L Zador. 1982. Asymptotic quantization error of continuous signals and the quantization dimension. <i>IEEE Transactions on Information Theory</i> , 28(2):139–149.	701
		702
		703
		704
	Chi Zhang, Huaping Zhong, Kuan Zhang, Chengliang Chai, Rui Wang, Xinlin Zhuang, Tianyi Bai, Qiu Jiantao, Lei Cao, Ju Fan, Ye Yuan, Guoren Wang, and Conghui He. 2025a. <a href="#">Harnessing diversity for important data selection in pretraining large language models</a> . In <i>The Thirteenth International Conference on Learning Representations</i> .	705
		706
		707
		708
		709
		710
		711
	Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, and 1 others. 2025b. Qwen3 embedding: Advancing text embedding and reranking through foundation models. <i>arXiv preprint arXiv:2506.05176</i> .	712
		713
		714
		715
		716
		717

Xinlin Zhuang, Jiahui Peng, Ren Ma, Yinfan Wang, Tianyi Bai, Xingjian Wei, Qiu Jiantao, Chi Zhang, Ying Qian, and Conghui He. 2025. Meta-rater: A multi-dimensional data selection method for pre-training language models. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10856–10896.

## A Unsupervised Hyperparameter Derivation

In Eq. (1), the scoring function relies on a weight vector  $\mathbf{w}$  to balance diverse geometric features. Instead of relying on heuristic grid search, we derive  $\mathbf{w}$  intrinsically from the data geometry. This process is grounded in the statistical stability analysis and spectral consensus visualized in Figure 7 and Figure 6.

### 1. Feature Stabilization (Log-Normal Priors).

Let  $\mathcal{F} = \{z_{\text{coh}}, z_{\text{len}}, z_{\text{ent}}, z_{\text{size}}\}$  be the set of raw feature vectors. As visualized in Figure 7, raw extensive metrics (Length, Size) exhibit heavy-tailed instabilities. To mitigate this, we first apply log transform to stabilize their magnitudes. We define the **stabilized feature matrix**  $\mathbf{Z}^\dagger$  as follows:

$$\begin{aligned} z_{\text{coh}}^\dagger &= z_{\text{coh}} \\ z_{\text{ent}}^\dagger &= z_{\text{ent}} \\ z_{\text{len}}^\dagger &= \log(z_{\text{len}}) \\ z_{\text{size}}^\dagger &= \log(z_{\text{size}}) \end{aligned} \quad (12)$$

Note that at this stage,  $z^\dagger$  represents the raw physical properties (e.g., larger  $z_{\text{len}}^\dagger$  still means longer sequence).

### 2. Standardization and Main Text Notation.

We then apply Z-score standardization to the entire matrix  $\mathbf{Z}^\dagger$  to ensure all dimensions share a unified scale (zero mean, unit variance). This yields the normalized metrics  $\tilde{\mathbf{z}}$  utilized in the main text (Eq. 1):

$$\tilde{z}_{k,f} = \text{Z-Score}(z_{k,f}^\dagger) \quad (13)$$

Thus,  $\tilde{\mathbf{z}}$  preserves the original polarity of the features. This is why Eq. 1 explicitly subtracts the penalty terms (Length, Entropy, Size) to convert them into a quality score.

**3. Spectral Weight Derivation.** To derive the consensus weights  $\mathbf{w}$ , we construct a temporary **Aligned Matrix**  $\mathbf{X}_{\text{align}}$  where all features are oriented towards "quality" (flipping the signs of terms):

$$\mathbf{X}_{\text{align}} = [\tilde{z}_{\text{coh}}, -\tilde{z}_{\text{ent}}, -\tilde{z}_{\text{len}}, -\tilde{z}_{\text{size}}] \quad (14)$$

(a) Geometric Consensus ( $\Sigma$ )

$z_{\text{len}}^\dagger$	1.02	0.40	0.20	0.33
$z_{\text{ent}}^\dagger$	0.40	1.02	0.60	0.40
$z_{\text{coh}}^\dagger$	0.20	0.60	1.02	0.28
$z_{\text{size}}^\dagger$	0.33	0.40	0.28	1.02
	$z_{\text{len}}^\dagger$	$z_{\text{ent}}^\dagger$	$z_{\text{coh}}^\dagger$	$z_{\text{size}}^\dagger$

(b) Derived Spectral Weights

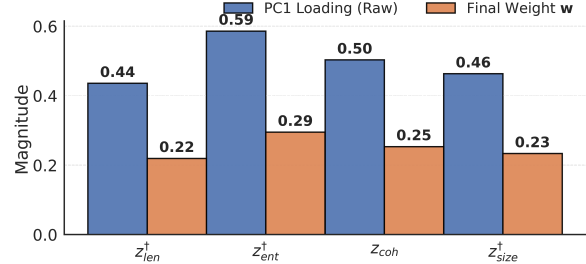


Figure 6: **Unsupervised Spectral Weight Derivation.** (a) The geometric consensus matrix  $\Sigma$ , computed on the **polarity-aligned metrics** (where penalties are inverted), reveals strong positive correlations (e.g., 0.60 between aligned entropy and cohesion). This suggests that these proxies share a common latent direction. (b) Instead of heuristic tuning, we derive the balancing coefficients  $\mathbf{w}$  directly from the first principal component (PC1). The spectral analysis intrinsically yields data-driven weights in which **Entropy** and **Cohesion** receive the largest coefficients.

We then compute the covariance matrix  $\Sigma = \frac{1}{K-1} \mathbf{X}_{\text{align}}^\top \mathbf{X}_{\text{align}}$ . As shown in Figure 6(a), this alignment reveals strong positive correlations. The first principal eigenvector  $\mathbf{v}_1$  of  $\Sigma$  captures the direction of **Maximum Consensus**. The final weights are derived by  $L_1$ -normalizing this vector:  $\mathbf{w} = \mathbf{v}_1 / \|\mathbf{v}_1\|_1$ .

## B Intrinsic Resolution Selection Details

### B.1 Metric Definition: Rank Stability

To robustly quantify the topological stability  $J_{\text{stab}}(K)$  (rank stability across resolutions), we utilize Kendall's Rank Correlation Coefficient ( $\tau$ ). This metric evaluates whether the relative quality ranking of clusters remains consistent after projecting to a finer resolution.

Given the two scoring vectors  $\mathbf{s}^{(K)}$  and  $\hat{\mathbf{s}}^{(K)}$  (from Eq. 3), we examine all possible pairs of clusters  $(i, j)$  where  $1 \leq i < j \leq K$ . A pair is classified based on the consistency of their relative

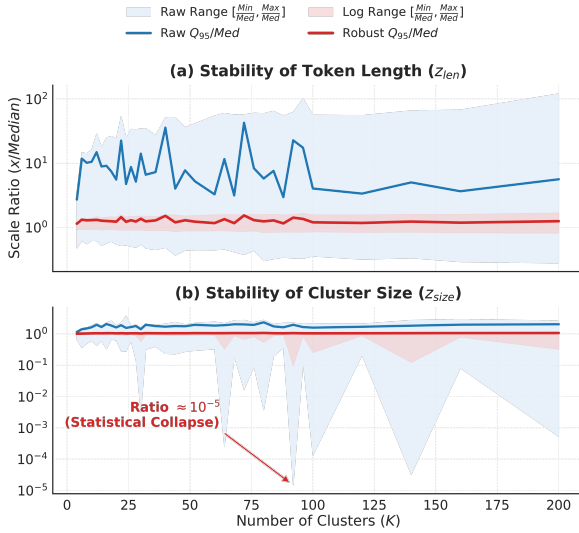


Figure 7: **Statistical Stability and the Necessity of Log Transform.** We analyze the stability of feature statistics across varying cluster resolutions ( $K \in [5, 200]$ ). (a) For token length, the raw scale (blue) exhibits high-frequency oscillations, while our **log-stabilized estimator** ( $z_{len}^\dagger$ , red) remains smooth. (b) For cluster size, the raw statistics suffer from extreme **statistical collapse** (highlighted by the arrow, where the **minimum-to-median ratio** drops to  $\sim 10^{-5}$ ), which would cause numerical instability in standard Z-score calculations.

ordering:

$$\text{Concordant} : (s_i - s_j)(\hat{s}_i - \hat{s}_j) > 0 \quad (15)$$

$$\text{Discordant} : (s_i - s_j)(\hat{s}_i - \hat{s}_j) < 0 \quad (16)$$

Let  $N_{conc}$  and  $N_{disc}$  denote the total counts of such pairs:

$$N_{conc} = \sum_{i < j} \mathbb{I}(\text{Concordant})$$

$$N_{disc} = \sum_{i < j} \mathbb{I}(\text{Discordant})$$

where  $\mathbb{I}(\cdot)$  is the indicator function. The final stability metric is the normalized difference between these counts:

$$J_{stab}(K) = \frac{N_{conc} - N_{disc}}{K(K-1)/2} = \frac{2(N_{conc} - N_{disc})}{K(K-1)} \quad (17)$$

The denominator  $K(K-1)/2$  represents the total number of unique pairs. Thus,  $J_{stab} \approx 1$  indicates that the geometric structure is perfectly preserved across resolutions.

## B.2 Implementation Optimizations

To further ensure stability against sample size variations, we implement the following engineering optimizations:

**1. Multi-Scale Hop Averaging.** Instead of relying solely on the immediate neighbor ( $K \rightarrow K+1$ ), which may be noisy, we compute stability across multiple strides  $\Delta K \in \{2, 4, 6\}$ . The final stability is a weighted average:

$$J_{final}(K) = \frac{\sum_{h \in \Delta K} \gamma_h \cdot J_{stab}(K \rightarrow K+h)}{\sum \gamma_h} \quad (18)$$

where  $\gamma_d$  are decay weights (e.g.,  $[0.5, 0.3, 0.2]$ ) that prioritize local consistency.

**2. Small-Sample Fisher Shrinkage.** When  $K$  is small, rank-correlation estimates can have high variance and appear overly optimistic. We therefore apply an **atanh-based shrinkage heuristic** to damp inflated stability scores in low- $K$  regimes. We treat  $J_{stab} \in (-1, 1)$  as a bounded rank-stability score and map it with  $z = \text{arctanh}(J_{stab})$  for shrinkage; this is an engineering correction rather than a statistical guarantee.

$$z_{shrunk} = z \cdot \tanh\left(\lambda_{shrink} \cdot \sqrt{N_{valid} - 3}\right), \quad (19)$$

$$r_{shrunk} = \tanh(z_{shrunk}). \quad (20)$$

where  $N_{valid}$  denotes the effective number of clusters participating in the rank comparison and  $\lambda_{shrink}$  is a regularization parameter.

## C Theoretical Proofs

In this section, we provide the detailed derivation for Theorem 1. We model the data selection process as an optimal quantization problem on a Riemannian manifold, drawing connections to recent theoretical advances in data pruning (Sorscher et al., 2022). All bounds below hold for an arbitrary number of clusters  $K$ ; in UniGeM we instantiate  $K = K^*$ , where  $K^*$  is selected by maximizing the stability objective  $J_{stab}(K)$  (Algorithm 1).

### C.1 Proof of Error Decomposition

Let  $Q : \mathcal{M} \rightarrow \{\mathbf{c}_1, \dots, \mathbf{c}_K\}$  be the quantization operator mapping any point  $x$  to its nearest global centroid  $\mathbf{c}_{k(x)}$ . Let  $C_k$  denote the Voronoi cell induced by  $\mathbf{c}_k$ , and define  $\alpha_k \triangleq \mu(C_k)$ .

Starting from the definition of the squared Wasserstein-2 distance (Peyré et al., 2019),

$$\begin{aligned} \mathcal{E}(S) &= W_2^2(\mu, \hat{\mu}_S) \\ &= \inf_{\gamma \in \Pi(\mu, \hat{\mu}_S)} \int_{\mathcal{M} \times \mathcal{M}} \|x - y\|^2 d\gamma(x, y) \end{aligned} \quad (21)$$

## Current Resolution $K$ Neighbor Resolution $K'$

Probabilistic Bridge  $\pi_{K \rightarrow K'}$

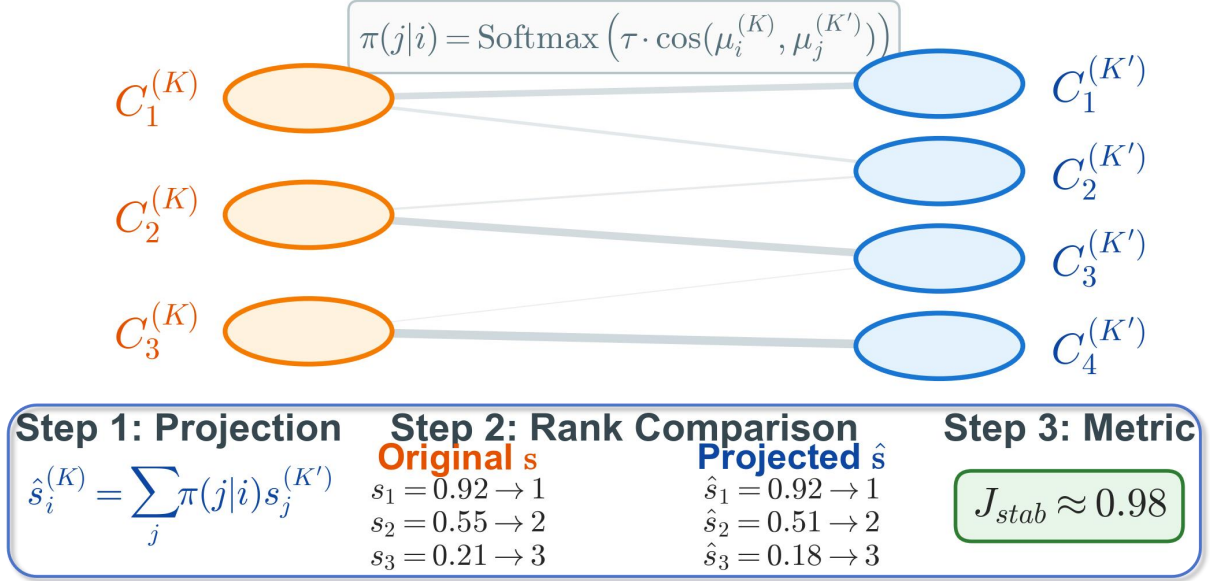


Figure 8: **Schematic of Intrinsic Resolution Selection.** The workflow proceeds in three stages corresponding to the mathematical derivation: (1) **Projection:** A probabilistic bridge  $\pi$  is constructed via centroid similarity (Eq. 2) to back-project scores from neighbor resolution  $K'$ , yielding the reconstructed profile  $\hat{s}^{(K)}$  (Eq. 3); (2) **Rank Comparison:** The relative ordering of the original geometric scores  $s^{(K)}$  is compared against the reconstructed proxies  $\hat{s}^{(K)}$ ; (3) **Metric Computation:** The final stability score  $J_{stab}$  is derived using Kendall's  $\tau$  rank correlation (Eq. 4). Numerical values (e.g., 0.92) are schematic examples for illustration.

we upper bound  $\mathcal{E}(S)$  by constructing an explicit two-stage transport plan through a centroid-supported intermediate measure. Define  $\mu_K \triangleq \sum_{k=1}^K \alpha_k \delta_{\mathbf{c}_k}$ . By the triangle inequality of  $W_2$ ,

$$W_2(\mu, \hat{\mu}_S) \leq W_2(\mu, \mu_K) + W_2(\mu_K, \hat{\mu}_S), \quad (22)$$

and using  $(a+b)^2 \leq 2a^2 + 2b^2$  yields the constant-factor bound

$$\mathcal{E}(S) \leq 2W_2^2(\mu, \mu_K) + 2W_2^2(\mu_K, \hat{\mu}_S). \quad (23)$$

**Stage-I (Global distortion).** Couple each  $x \sim \mu$  with its centroid  $Q(x)$  to obtain

$$W_2^2(\mu, \mu_K) \leq \sum_{k=1}^K \int_{C_k} \|x - \mathbf{c}_k\|^2 d\mu(x) \triangleq \mathcal{T}_1. \quad (24)$$

**Stage-II (Within-cluster residual).** For analysis, we view the selected-set empirical measure as a cluster-wise mixture

$$\hat{\mu}_S \triangleq \sum_{k=1}^K \alpha_k \hat{\mu}_{S_k}, \quad (25)$$

where  $\hat{\mu}_{S_k}$  is the empirical measure supported on  $S_k$ . Let  $S_k \triangleq S \cap C_k$  be the selected subset inside

cluster  $C_k$ . The second term  $W_2^2(\mu_K, \hat{\mu}_S)$  measures how well the selected points within each  $C_k$  represent the local mass anchored at  $\mathbf{c}_k$ . This induces a within-cluster residual term that we summarize by

$$W_2^2(\mu_K, \hat{\mu}_S) \leq \sum_{k=1}^K \alpha_k \mathbb{V}(S_k) \triangleq \mathcal{T}_2. \quad (26)$$

where we define the within-cluster residual energy (a second moment w.r.t. the centroid) as

$$\mathbb{V}(S_k) \triangleq \mathbb{E}_{x \sim \hat{\mu}_{S_k}} \|x - \mathbf{c}_k\|^2. \quad (27)$$

Combining the two parts with Eq. (23) yields a constructive decomposition of  $\mathcal{E}(S)$  into a global quantization term  $\mathcal{T}_1$  and a local within-cluster term  $\mathcal{T}_2$ , up to constant factors commonly used in quantization-style analyses (Gray and Neuhoff, 2002).

**Remark (High-dimensional intuition).** In high-dimensional embeddings ( $d \gg 1$ ), cross-terms between centroid error and within-cluster residual are often empirically small due to concentration effects (Ledoux, 2001), which motivates the near-additive behavior observed in practice; however, our bound above does not rely on this approximation.

## 880 C.2 Bound Derivation for Stage-I ( $\mathcal{T}_1$ )

881 The first term  $\mathcal{T}_1$  corresponds to the classical  
882 high-resolution quantization error. According to  
883 **Zador’s Theorem** (Zador, 1982), for a quantizer  
884 with  $K$  codepoints on a  $d$ -dimensional manifold  
885 with probability density function  $p(x)$ , the asymp-  
886 totic distortion satisfies:

$$887 \lim_{K \rightarrow \infty} K^{2/d} \cdot \mathcal{T}_1 = J_d \|p\|_{d/(d+2)}, \quad (28)$$

888 where  $J_d$  is the coefficient of the optimal lattice  
889 quantizer in  $\mathbb{R}^d$ . This yields the Stage-I term in our  
890 bound:

$$891 \mathcal{T}_1 \leq C_d \cdot K^{-2/d}. \quad (29)$$

892 **Implication:** This suggests that Stage-I controls  
893 global covering distortion through the choice of  $K$ .

## 894 C.3 Bound Derivation for Stage-II ( $\mathcal{T}_2$ )

895 The second term  $\mathcal{T}_2$  represents the intra-cluster  
896 residual variance. For a standard random sam-  
897 pler, this corresponds to the raw cluster variance,  
898 which we denote by  $\sigma_k^2$  for cluster  $C_k$ . Let  $p_k(x) \triangleq$   
899  $p(x)/\alpha_k$  for  $x \in C_k$  denote the conditional density  
900 within cluster  $C_k$ , where  $\alpha_k = \mu(C_k)$ .

901 **Assumption (dominant structural filtering).**  
902 Within each cluster  $C_k$ , we assume the auxiliary  
903 reweighting terms used in practice (e.g., probe  
904 score  $P_{S_j}$  and cohesion gate  $\beta_{S_j}$ ) are either (i) ap-  
905 proximately independent of the radial deviation  
906  $\|x - \mathbf{c}_k\|$  or (ii) bounded and do not systematically  
907 favor higher-deviation points. Under this assump-  
908 tion, the dominant geometric effect of Stage-II is  
909 governed by  $\exp(-\lambda \mathcal{L}_{struct}(x))$ .

910 UniGeM modulates the sampling probability via  
911  $P(x) \propto \exp(-\lambda \mathcal{L}_{struct}(x))$ . For theoretical anal-  
912 ysis, we approximate this soft exponential decay as  
913 a truncation mechanism on an effective acceptance  
914 region

$$915 \Omega_{UniGeM} = \{x \in C_k \mid \mathcal{L}_{struct}(x) < \tau\}, \quad (30)$$

916 where  $\tau$  is a confidence threshold implicitly con-  
917 trolled by  $\lambda$ .

918 Define the random-baseline within-cluster sec-  
919 ond moment as

$$920 \sigma_k^2 \triangleq \int_{C_k} \|x - \mathbf{c}_k\|^2 p_k(x) dx. \quad (31)$$

921 Under the truncation approximation, UniGeM  
922 induces the conditional density  $q_k(x) \triangleq$

$\frac{p_k(x) \mathbf{1}_{[x \in \Omega_{UniGeM}]}}{Z_k}$  with  $Z_k \triangleq \int_{\Omega_{UniGeM}} p_k(x) dx$ ,  
and the resulting second moment is

$$923 \mathbb{V}(S_{UniGeM}^{(k)}) \approx \int_{\Omega_{UniGeM}} \|x - \mathbf{c}_k\|^2 q_k(x) dx. \quad (32)$$

924 We define the pruning gain as

$$925 \Delta_{gain}^{(k)} \triangleq \sigma_k^2 - \int_{\Omega_{UniGeM}} \|x - \mathbf{c}_k\|^2 q_k(x) dx \geq 0, \quad (33)$$

926 where the non-negativity holds when the accep-  
927 tance region preferentially keeps lower-deviation  
928 points.

## 931 C.4 Final Theorem Assembly and Remark

932 Substituting the bounds for  $\mathcal{T}_1$  and  $\mathcal{T}_2$  back into the  
933 decomposition yields:

$$934 \mathcal{E}(S_{UniGeM}) \leq 2C_d K^{-2/d} + 2 \sum_{k=1}^K \alpha_k (\sigma_k^2 - \Delta_{gain}^{(k)}). \quad (34)$$

935 **Practical proxy.** In Eq. (11), the gain  $\Delta_{gain}^{(k)}$   
936 captures the variance mass removed by pruning  
937 within  $C_k$ . In practice, high-deviation samples cor-  
938 respond to large standardized feature deviations;  
939 thus  $\mathcal{L}_{struct}(x)$  (a rectified squared Mahalanobis  
940 distance) directly serves as a proxy for identify-  
941 ing the rejected region  $\Omega^c$  and monitoring pruning  
942 strength.

## 943 D Annotation Model.

944 We leverage the Qwen3-235B model as a **Knowl-  
945 edge Probe** to inspect the semantic attributes of  
946 these samples. Adopting a **Model-Based Anno-  
947 tation** strategy (Seed et al., 2025), we design a  
948 structured system prompt to decompose the anal-  
949 ysis into four complementary dimensions: *Code  
950 Quality, Engineering Design, Training Suitability,*  
951 *and Knowledge Density*. This multi-dimensional  
952 rubric aims to capture both syntactic correctness  
953 and training-relevant content.

954 The specific system prompt used for this annota-  
955 tion is provided below.

### System Prompt: Code Data Evaluation Strategy

You are an expert code evaluator. Your task is to assess the provided code snippet based on four distinct dimensions. For each dimension, assign a score on a scale of 1 to 5 according to the criteria below.

**1. Code Quality & Compliance** (Focus: Syntax, nam-

ing, and readability)

- **1-2 (Low):** Severe syntax errors or logical flaws; chaotic naming conventions; inconsistent indentation; contains dead code or empty functions; non-executable.
- **3-4 (Mid):** Syntactically correct and executable; reasonable naming and formatting; minor redundancies or loose structure but generally readable.
- **5 (High):** Error-free; strict adherence to standards (e.g., PEP8); precise naming; concise logic with no redundancy; highly readable and linear structure.

## 2. Algorithmic & Engineering Design (Focus: Modularity, robustness, and system thinking)

- **1-2 (Low):** Monolithic structure (global scope); lack of modularity; absence of input validation or exception handling; hardcoded constants; fails on edge cases.
- **3-4 (Mid):** Basic modularity (function splitting); separation of concerns; foundational error checking; clear structure but lacks advanced abstraction or extensibility.
- **5 (High):** High-level abstraction (classes/patterns); robust engineering (exception safety, resource management); extensible, testable, and demonstrates system-level thinking.

## 3. Training Suitability (Focus: Educational value and style consistency for LLMs)

- **1-2 (Low):** Arbitrary naming; missing or misleading comments; mixed styles; hallucinated or fragmentary code; lacks context.
- **3-4 (Mid):** Normative naming; comments cover key steps; consistent style; readable logic suitable for beginners or general training.
- **5 (High):** Strict adherence to community best practices; insightful comments explaining design intent; exemplary structure suitable for teaching or high-quality fine-tuning.

## 4. Knowledge Density (Focus: Technical insights and domain expertise)

- **1-2 (Low):** Trivial operations (e.g., simple I/O, loops); lacks domain knowledge; low information density equivalent to introductory tutorials.
- **3-4 (Mid):** Standard algorithms or patterns (e.g., DFS, HashMaps); implements common best practices; explicit but generic knowledge.
- **5 (High):** Non-trivial insights or cross-domain knowledge (e.g., memory alignment, bitwise optimization, lock-free structures); reveals underlying principles or deep optimizations.

**Output Format:** Output strictly in valid JSON format: {"code\_quality": int, "algorithm\_and\_engineering": int, "training\_suitability": int, "knowledge\_score": int}

**Data Input:** \$content

## E Experiment Details

In this appendix, we provide the condensed technical specifications for our experiments, including the model architecture, training recipe, and adaptation protocols for SOTA baselines.

### E.1 Model Architectures

We employ a fine-grained sparse MoE architecture with identical expert parameterization across mod-

els. We scale total capacity by increasing the number of experts, while keeping the routing strategy fixed (Top-2) so that the per-token activated parameter budget remains constant, preserving inference throughput.

Configuration	UniGeM-8B	UniGeM-16B
Total / Active Params	8.0B / 1.4B	16.8B / 1.4B
Total Experts ( $N$ )	32	64
Routing Strategy	Top-2	Top-2
Hidden / Layers	2048 / 24	2048 / 24

### E.2 Pre-training Configuration

The models are trained from scratch on a 100B-token mixture, consisting of 70B code tokens and 30B code-related text tokens. We utilize **Qwen3-235B** to retrieve the text component from Common Crawl to ensure semantic relevance. The training employs a Warmup-Stable-Decay (WSD) schedule.

**Hyperparameter Selection Logic.** The geometric hyperparameters  $\{\lambda, T_{scale}, \epsilon\}$  are calibrated based on the statistical moments of the 20% probe manifold. We set the **structural**  $\lambda = 0.5$  to ensure that samples deviating beyond  $2\sigma$  from the geometric consensus incur a significant weight reduction (e.g. reduced to  $\approx 13.5\%$  of the original weight), effectively pruning logical outliers. The **scale factor**  $T_{scale} = 20$  is employed to amplify the density contrast during cross-resolution alignment, ensuring the stability-driven clustering captures sharp manifold boundaries. Finally, an **exploration floor**  $\epsilon = 0.01$  is maintained to preserve long-tail distributional diversity and mitigate manifold approximation errors.

Hyper-parameter	Value
Optimizer	AdamW ( $\beta_1 = 0.9, \beta_2 = 0.95, WD=0.1$ )
Peak Learning Rate	$3.0 \times 10^{-4}$ (8B) / $2.4 \times 10^{-4}$ (16B)
Batch Size / Seq Len	2,560 $\rightarrow$ 8,960 / 4,096 Tokens
Stability Mechanisms	NormHead, Stochastic Routing Warmup

UniGeM Geometric Mining Params	
Structural ( $\lambda$ )	0.5 (calibrated to $2\sigma$ consensus)
Exploration Floor ( $\epsilon$ )	0.01 (1% diversity reserve)
Transition Scale ( $T_{scale}$ )	20 (Eq. 2)
Clustering & Embedding	
Embedding Model	<i>Qwen3-embedding</i>
Global Resolution ( $K^*$ )	72 (via stability-driven selection)
K-means Iterations	10 (Stage-I) / 5 (Stage-II)

### E.3 SOTA Baseline Adaptation

Standard implementations of CLIMB and Meta-rater are designed for general text and often rely on proxy tasks (e.g., MMLU) or generic encoders (e.g., BERT), which can introduce domain mismatch for code. Using these methods strictly off-the-shelf would make the comparison less informative. We therefore apply a **Code Adaptation** protocol to improve their alignment with code and obtain strong code-aware baselines.

- **Nemotron-CLIMB Adaptation:** We replace general BERT embeddings with **Qwen3-embedding** to perform code-aware semantic clustering ( $K = 100$ ). The optimization target is shifted from MMLU to a **Code Oracle** ( $V_{\text{code}}$ ), which computes weighted validation loss on **MBPP-Sanitized**, **HumanEval-Pack**, and logic-dense samples from **DS-1000**.
- **Meta-rater Adaptation:** We re-define the PRRC framework into **Code-PRRC**, focusing on *Professionalism* (complexity), *Readability* (style), *Reasoning* (flow density), and *Cleanliness* (syntax). We score 500k seed samples via **Qwen3-235B** to distill four specialized **ModernBERT-base** scorers capable of handling long-context code quality assessment.
- **Common Proxy Setup:** To search for optimal weights, both methods utilize a **350M Dense Transformer** as a proxy model. These proxies are trained on 2B token slices across multiple iterations (64 trials for CLIMB; 256 for Meta-rater) to fit a LightGBM-based quality-to-loss regressor.

### F Distributional Analysis of Features for Stage-I Cluster

In this section, we provide a detailed distributional analysis of the geometric proxies observed in our large-scale pre-training experiment. Specifically, we visualize the feature statistics across the optimal resolution of  $K^* = 72$  **global clusters**, identified via the Topological Stability analysis in Section 4.3. These empirical results validate the statistical assumptions underpinning our Geometric Scoring function (Eq. 1).

**Transformation of Extensive Properties.** As illustrated in the top row of Figure 9, the raw distributions of **Cluster Size** ( $z_{\text{size}}$ ) and **Sequence Length** ( $z_{\text{len}}$ ) exhibit significant right-skewness across the 72 latent domains. This confirms our hypothesis in

Section 2.2 that the raw code corpus is highly heterogeneous, spanning multiple orders of magnitude. Direct usage of these raw metrics would result in *variance dominance*, where spectral analysis is biased by magnitude outliers rather than structural quality. However, as shown in the bottom row, applying the logarithmic transformation effectively projects these features into log-space, where the transformed values are closer to Gaussian. This supports the normality assumption behind Z-score standardization and the subsequent spectral consensus step.

**Stability of Intensive Properties.** Conversely, the intensive properties—**Cohesion** ( $z_{\text{coh}}$ ) and **Entropy** ( $z_{\text{ent}}$ )—naturally display bounded, unimodal distributions across the 72 clusters without transformation. This implies that the experimentally identified clusters are statistically well-formed, maintaining consistent internal densities and semantic purities. Consequently, we validate our hybrid processing strategy: while extensive features require logarithmic dampening to mitigate scale disparities, intensive features can be directly utilized as linear geometric priors to preserve their original sensitivity.

### G Distributional Analysis of Semantic Scoring

In this section, we report the statistical characteristics of the semantic scores obtained during the Stage-II exploration. We visualize both the raw discrete ratings from the Annotation Model and the aggregated continuous scores for sub-clusters to provide a comprehensive view of the data quality distribution.

#### G.1 Probe and Sub-Cluster Score Distributions

Figure 10 presents the distribution of discrete quality ratings (scale 1–5) assigned by the Annotation Model to individual probe samples. The results indicate a pronounced left-skewed distribution across all four evaluation dimensions. Specifically, the *Code Quality* and *Training Suitability* metrics show a high concentration of samples receiving perfect or near-perfect scores ( $> 60\%$  rated as 5).

Figure 11 further illustrates the probability density of the aggregated Semantic Score ( $P_{S_j}$ ) for sub-clusters. Consistent with the probe-level observations, the cluster-level scores are heavily concentrated in the high-value interval  $[4.0, 5.0]$ . This

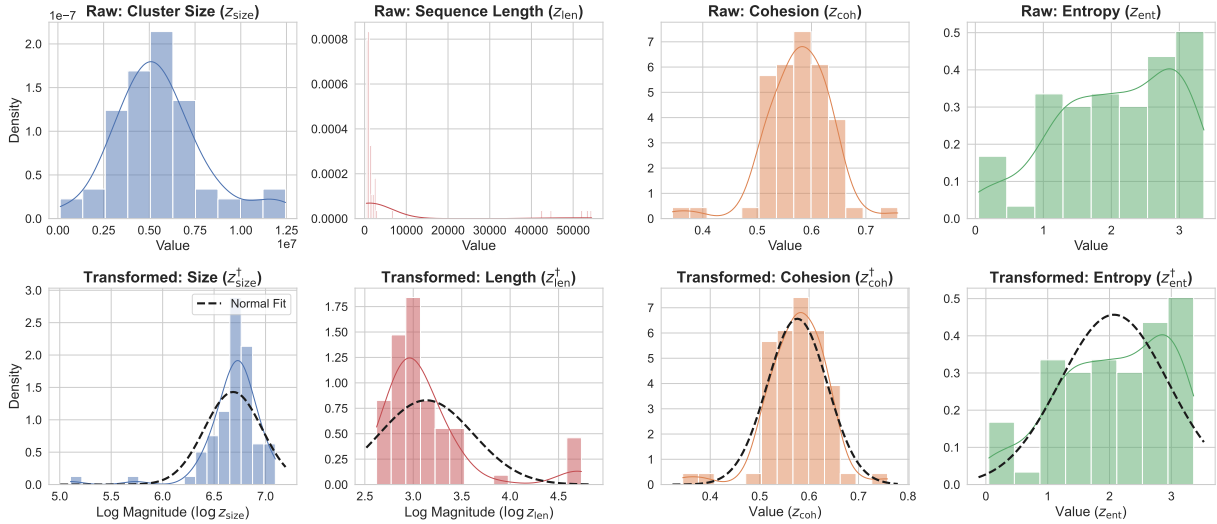


Figure 9: **Distributional Transformation across the  $K^* = 72$  Global Clusters.** (Top Row) The raw distributions of extensive properties—*Cluster Size* ( $z_{\text{size}}$ ) and *Sequence Length* ( $z_{\text{len}}$ )—exhibit extreme heavy-tailed skewness across the 72 experimental clusters. (Bottom Row) Applying the logarithmic transformation stabilizes these features; their log-values are closer to a Gaussian shape (as illustrated by the overlaid normal fits). In contrast, the intensive properties *Cohesion* and *Entropy* naturally follow a stable unimodal distribution, confirming the robustness of the extracted latent manifolds.

1091 suggests that the majority of the code corpus, after  
 1092 initial filtering, is perceived as syntactically valid  
 1093 and high-quality by the LLM judge.

## 1094 G.2 Implications for Selection Strategy

1095 The observed "ceiling effect" in the score distri-  
 1096 butions highlights a potential limitation in relying  
 1097 solely on semantic scoring: the lack of discrimina-  
 1098 tion in the high-score regime. Since a significant  
 1099 portion of sub-clusters achieves a saturated score  
 1100 ( $P_{S_j} \approx 5.0$ ), semantic metrics alone may strug-  
 1101 gle to differentiate between intrinsic high-value  
 1102 domains and superficially correct samples. This  
 1103 empirical observation motivates the design of the  
 1104 **Structural** ( $\mathcal{L}_{\text{struct}}$ ) within the UniGeM frame-  
 1105 work. By introducing geometric constraints as an  
 1106 orthogonal selection criterion, UniGeM effectively  
 1107 handles these saturated distributions, filtering out  
 1108 sub-clusters that appear high-quality but are topo-  
 1109 logically inconsistent with the domain manifold.

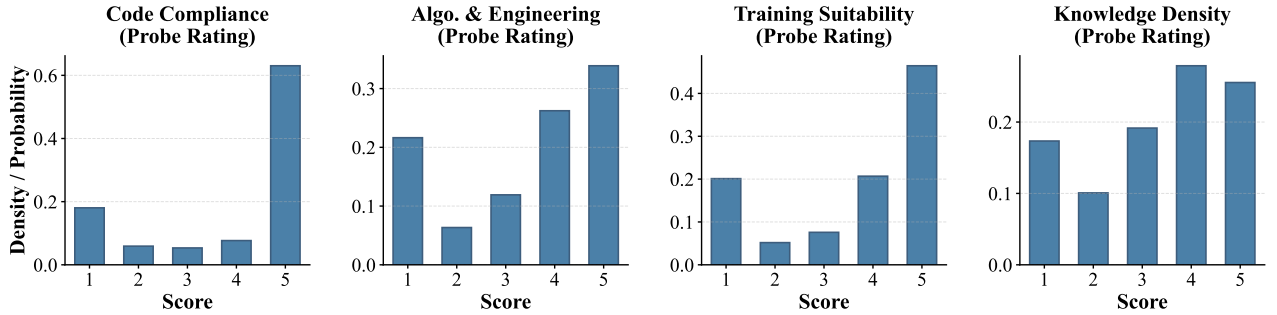


Figure 10: **Discrete Probe Quality Distribution (Raw Prompt Outputs)**. The probability mass distribution of integer scores (1–5) assigned by the **Annotation Model** to probe samples. The data exhibits a strong skew towards the upper bound (scores 4 and 5) across all dimensions, reflecting the general high acceptance rate of the judge model for code syntax.

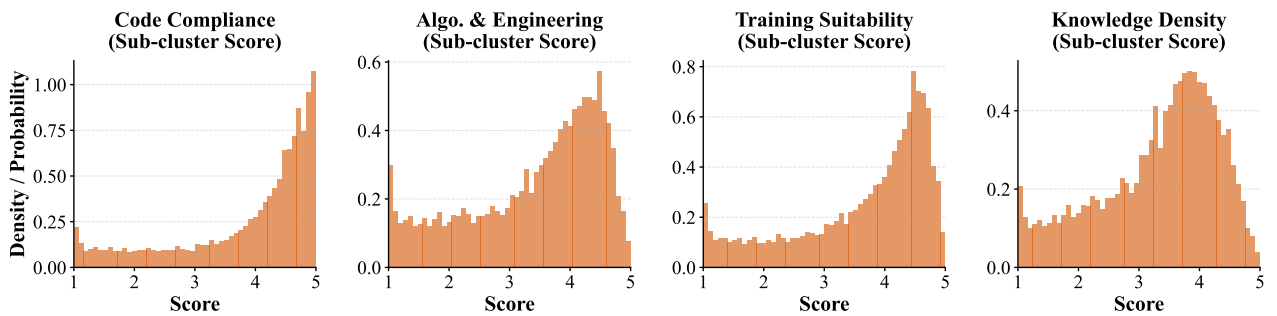


Figure 11: **Sub-cluster Semantic Score Distribution ( $P_{S_j}$ )**. Kernel density estimation of the aggregated semantic scores for sub-clusters  $S_j$ . The distribution shows a saturation effect in the 4.0 – 5.0 range. This lack of variance in the top quantile underscores the necessity of incorporating geometric penalties ( $\mathcal{L}_{struct}$ ) to introduce discriminative gradients among high-scoring clusters.

## 1110 H Reproducibility Statement

1111 We provide sufficient details for an independent  
1112 reimplementaion of UniGeM and for reproducing  
1113 our experimental results.

1114 **Method / Algorithm.** The full method specifi-  
1115 cation (problem setup, Stage-I Macro-Exploration,  
1116 Stage-II Micro-Mining, and all scoring terms) is  
1117 described in Section 2.

1118 **Experimental Protocol.** The evaluation pro-  
1119 tocol (corpus construction, sampling ratio, model  
1120 settings, and compared baselines) is summarized  
1121 in Section 3. Training and evaluation hyperparam-  
1122 eters, together with the calibration logic for the  
1123 geometric parameters, are provided in Appendix E.

1124 **Baselines and Prompts.** Details of the code  
1125 adaptaion for CLIMB and Meta-rater are docu-  
1126 mented in Appendix E. The probe/annotation rubric  
1127 and the exact system prompt used for semantic scor-  
1128 ing are provided in Appendix D.

1129 **Code Release.** We will release the complete  
1130 codebase (data processing, clustering/selection,  
1131 training scripts, and evaluation) in a public reposi-  
1132 tory upon completion of an internal review.

1133 **Compute and infrastructure.** All pre-training  
1134 runs were executed on 64 NVIDIA H800 GPUs.  
1135 Each 100B-token pre-training run took approx-  
1136 imately 15 hours wall-clock time. Automated  
1137 benchmark evaluation was executed on NVIDIA  
1138 H20 GPUs.

1139 **Use of AI assistants.** AI assistants were used  
1140 only for minor language polishing and clarity im-  
1141 provements during writing. All scientific content  
1142 was developed and verified by the authors.